

COMPUTATIONAL INTELLIGENCE

(INTRODUCTION TO MACHINE LEARNING) SS16

Lecture 2:

- Linear Regression
- Gradient Descent
- Non-linear basis functions

LINEAR REGRESSION MOTIVATION

Why Linear Regression?

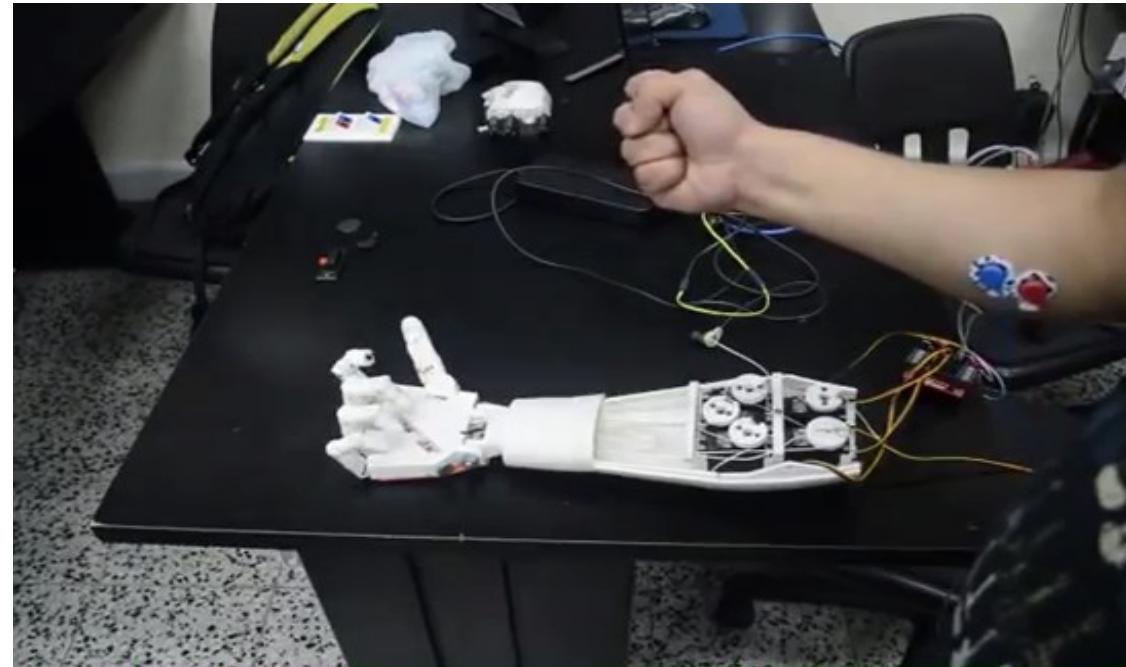
- Regression = Prediction **of real-valued outputs**
- Simplest regression algorithm
 - Easy, and fast
 - **Benchmark** algorithm
- **Mathematical Concepts** introduced
 - Data format and Matrix notation
 - Minimizing a **cost function**: gradient descent
 - Non-linear features and basis functions

Examples: (linear) regression application

- Social science: relationship between data
- Brain computer interfaces
- Neuroprosthetic control

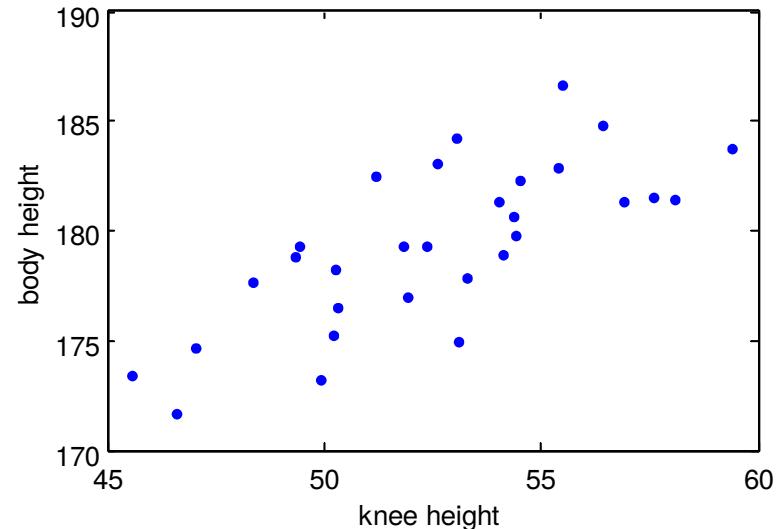
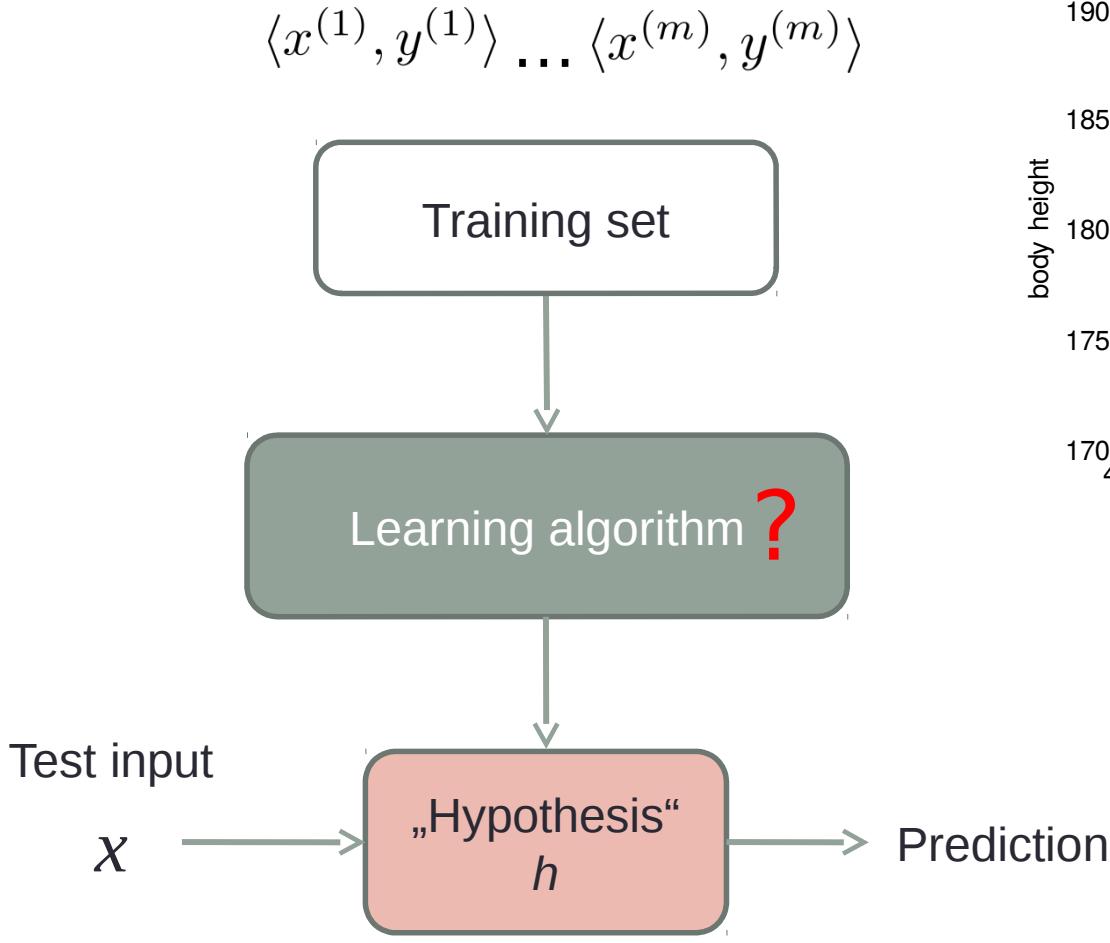
Examples: (linear) regression application

- Social science: relationship between data
- Brain computer interfaces
- Neuroprosthetic control



LINEAR REGRESSION WITH ONE INPUT

Linear regression with one input



Hypothesis

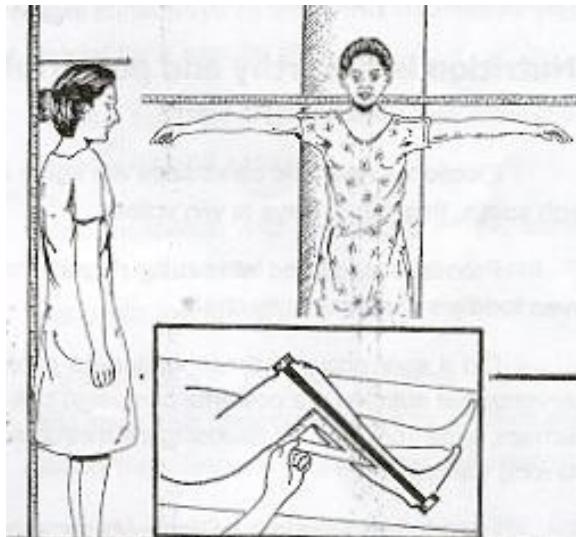
$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

Parameters ?

$$\theta = (\theta_0, \theta_1)$$

A regression problem

- We want to learn to predict a **person's height** based on his/her **knee height** and/or **arm span**
- This is useful for patients who are **bed bound** or in a wheelchair and cannot stand to take an accurate measurement of their height

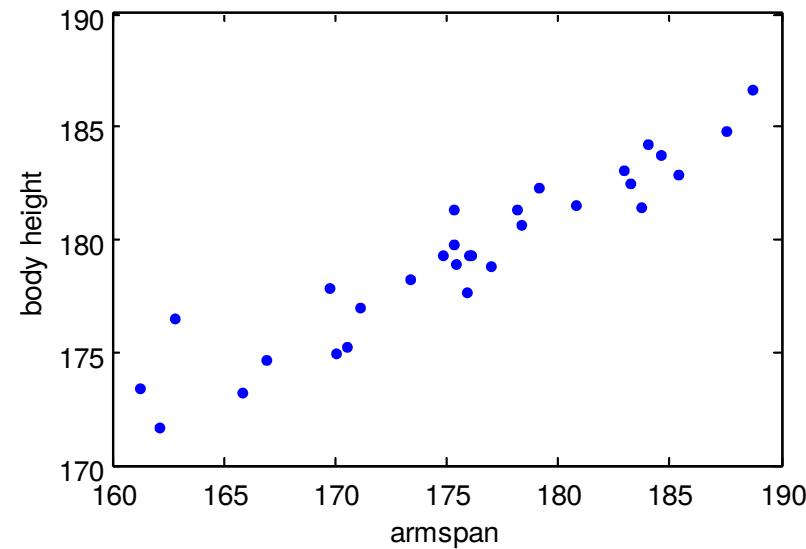
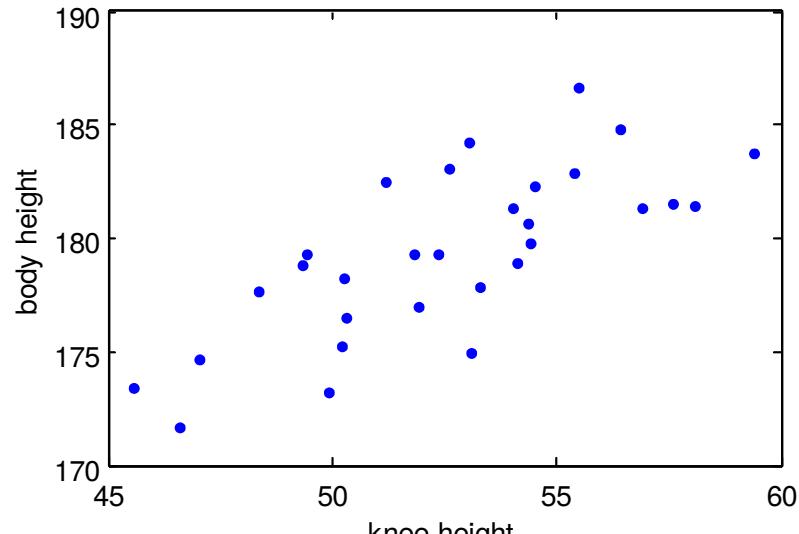


Knee Height [cm]	Arm span [cm]	Height [cm]
50	166	171
56	172	175
52	174	168
...

Example Data

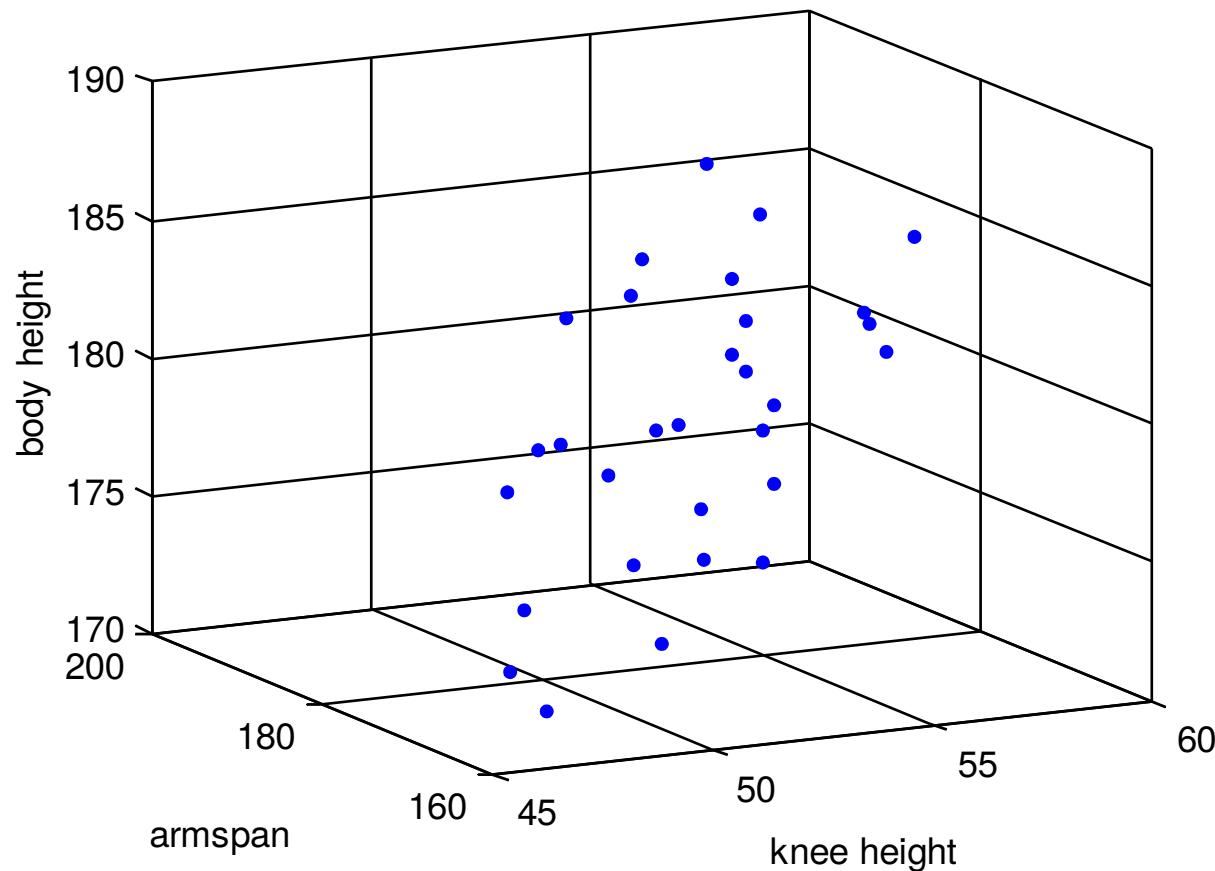
Knee height [cm]	Arm span [cm]	Height [cm]
50	166	171
56	172	175
52	174	168
...

$m=30$ data points



Example Data

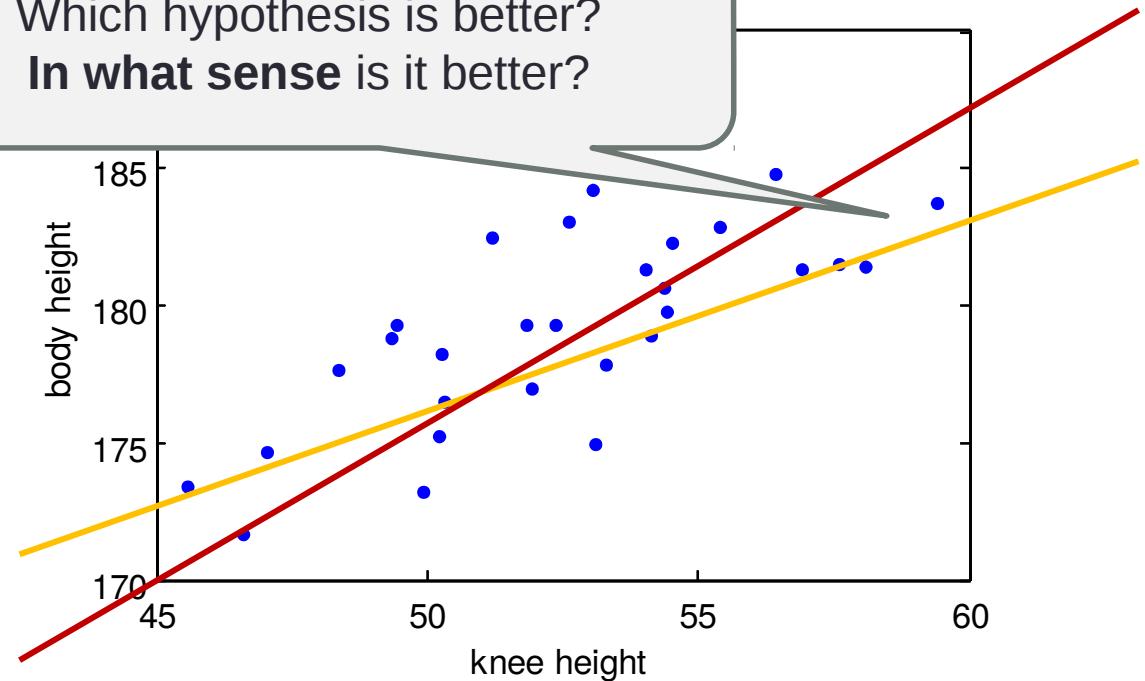
Knee Height [cm]	Arm span [cm]	Height [cm]
50	166	171
56	172	175
52	174	168
...



Linear regression with one input

Knee Height [cm]	Height [cm]
50	171
56	175
52	168
...	...

Which hypothesis is better?
In what sense is it better?



Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

Parameters ?

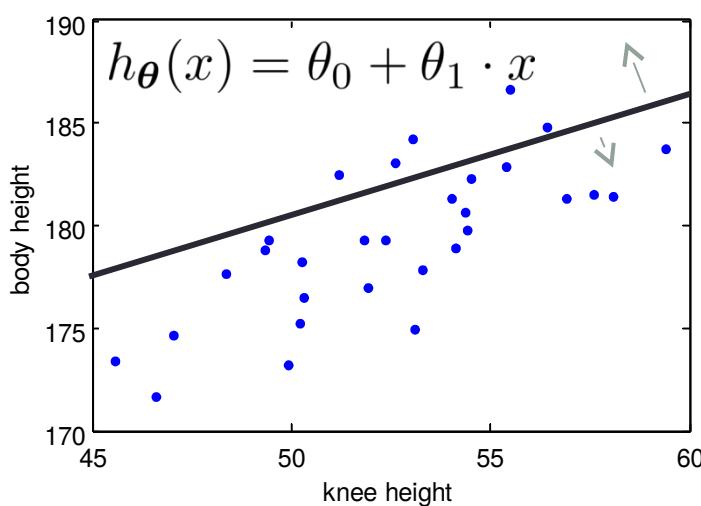
$$\theta = (\theta_0, \theta_1)$$

Formalization of problem

Knee Height [cm]	Height [cm]
50	171
56	175
52	168
...	...

$$x^{(i)} \quad y^{(i)}$$

m=30 data points



- Given m training examples

$$\langle x^{(1)}, y^{(1)} \rangle \dots \langle x^{(m)}, y^{(m)} \rangle$$

- Goal: learn parameters

$$\boldsymbol{\theta} = (\theta_0, \theta_1)$$

such that

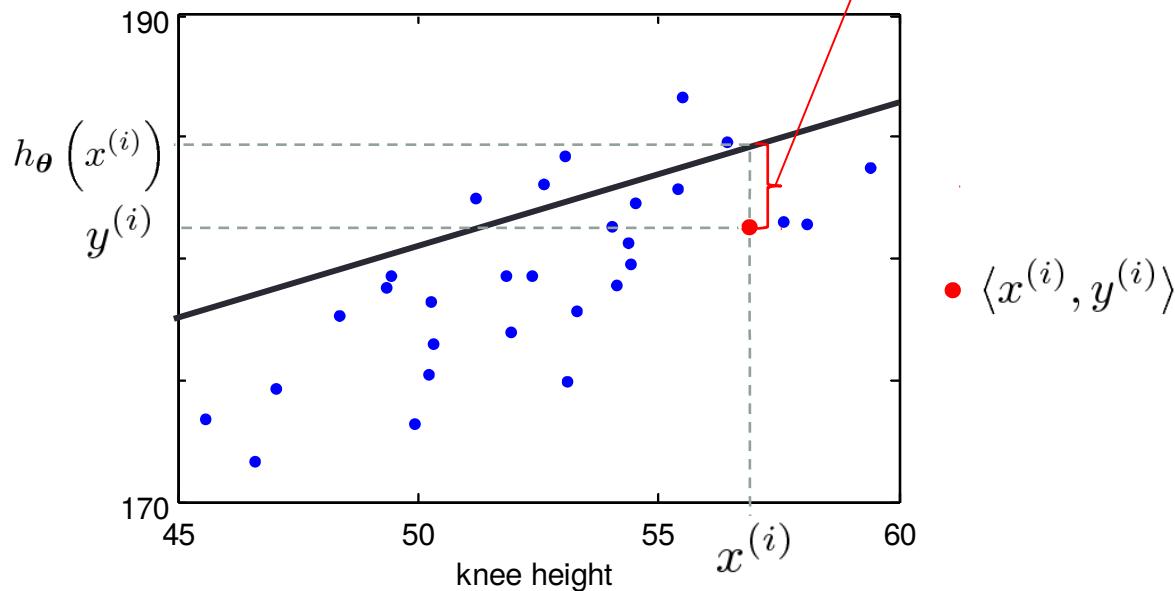
$$h_{\boldsymbol{\theta}}(x^{(i)}) = \theta_0 + \theta_1 \cdot x^{(i)} \approx y^{(i)}$$

for all training examples $i=1\dots 30$.

Least Squares Objective

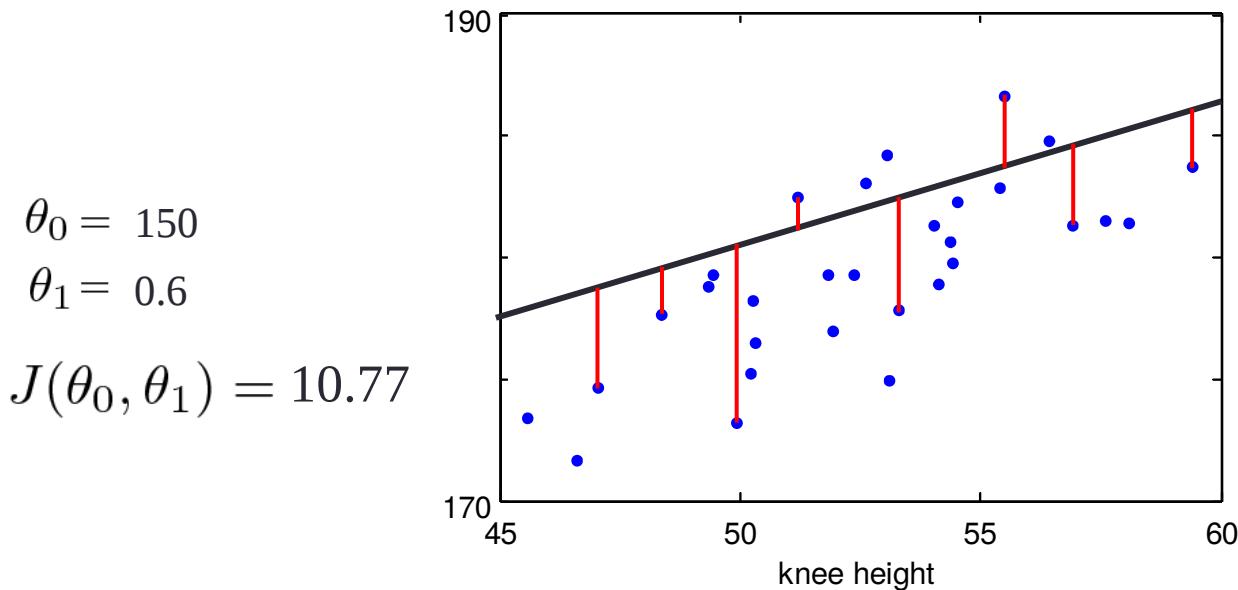
- Minimize Error $J(\theta_0, \theta_1) = \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$

$$\begin{aligned}\theta_0 &= 150 \\ \theta_1 &= 0.6\end{aligned}$$



Least Squares Objective

- Minimize Error
$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$
- cost function
- mean squared error

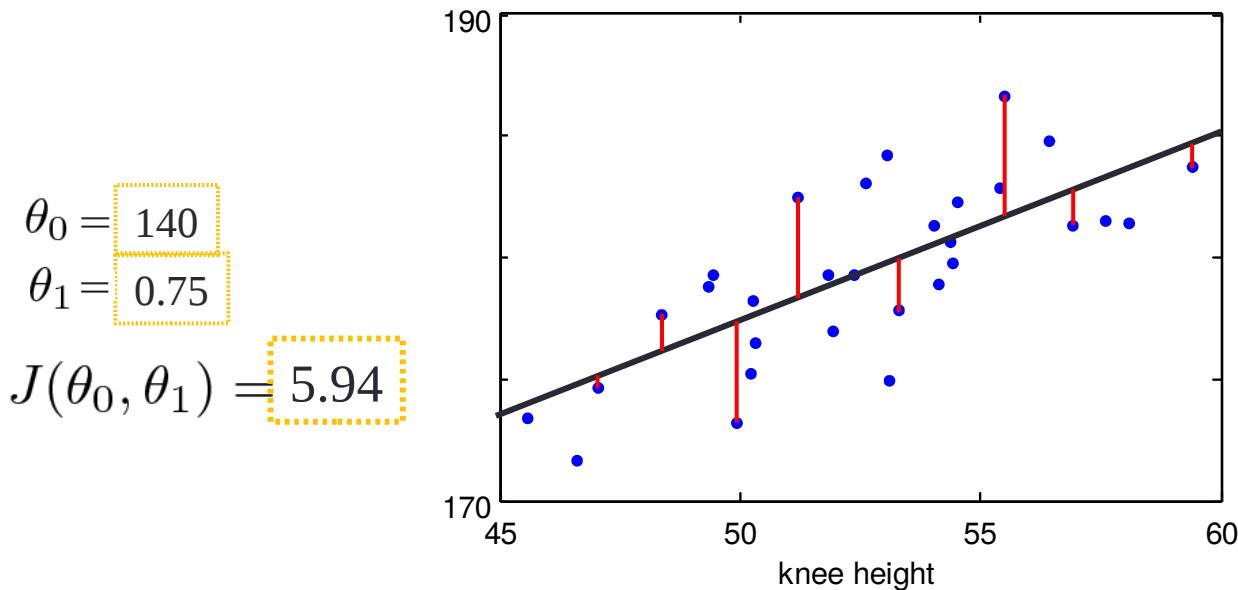


Least Squares Objective

- Minimize Error $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \underline{\left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2}$

cost function

mean squared error

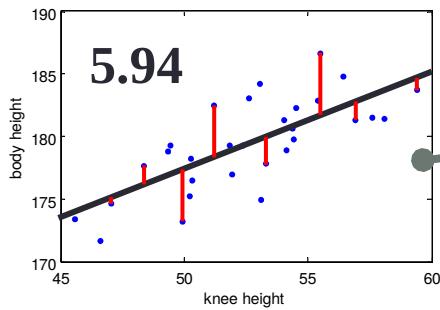


Cost function illustrated

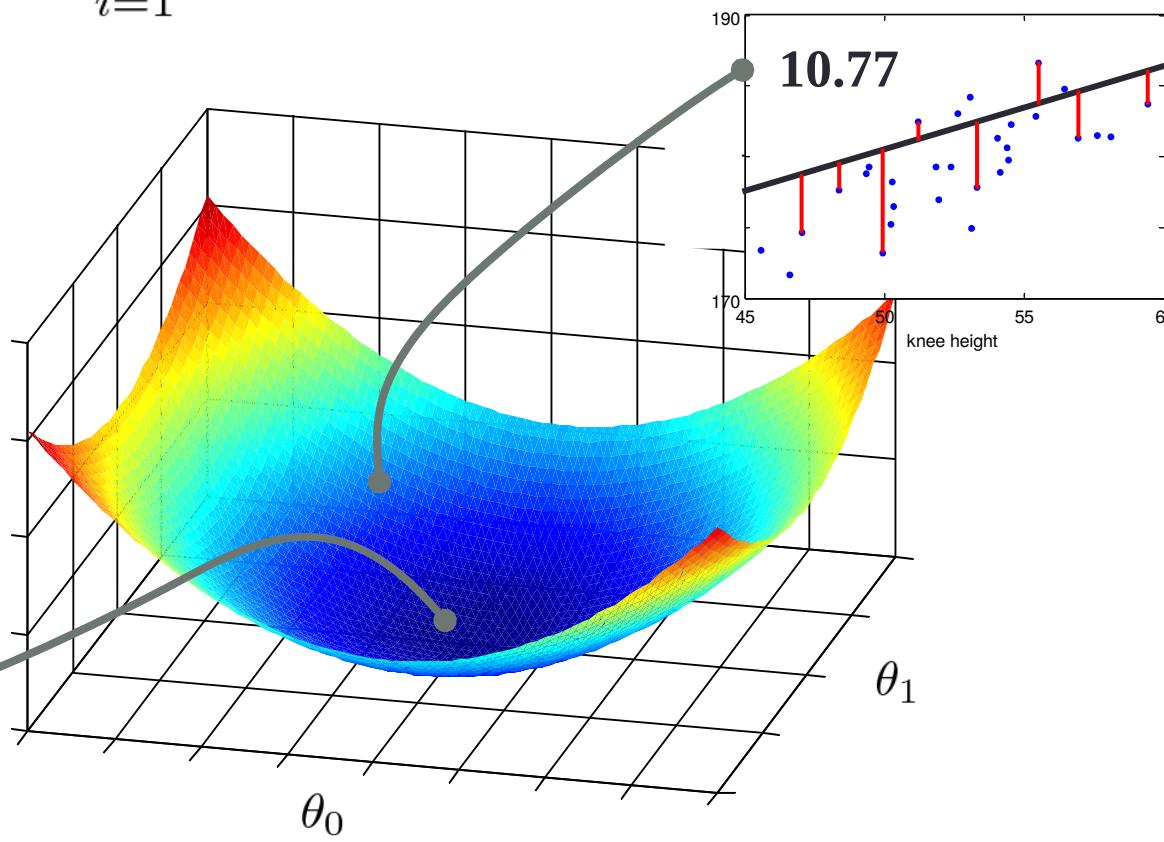
$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

Properties of cost function:

- **Quadratic** function
- **Convex** „Bowl“-shaped
- Unique local and global minimum (under „regular“ conditions)



$$J(\theta_0, \theta_1)$$



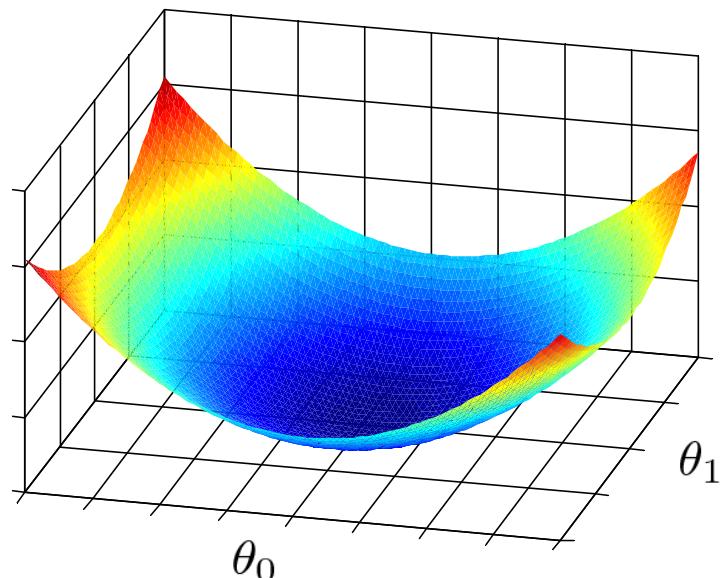
Minimizing the cost

- Two ways to find the parameters $\theta = (\theta_0, \theta_1)$ minimizing

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

- Gradient descent
- Direct analytical solution
(setting derivatives = 0)

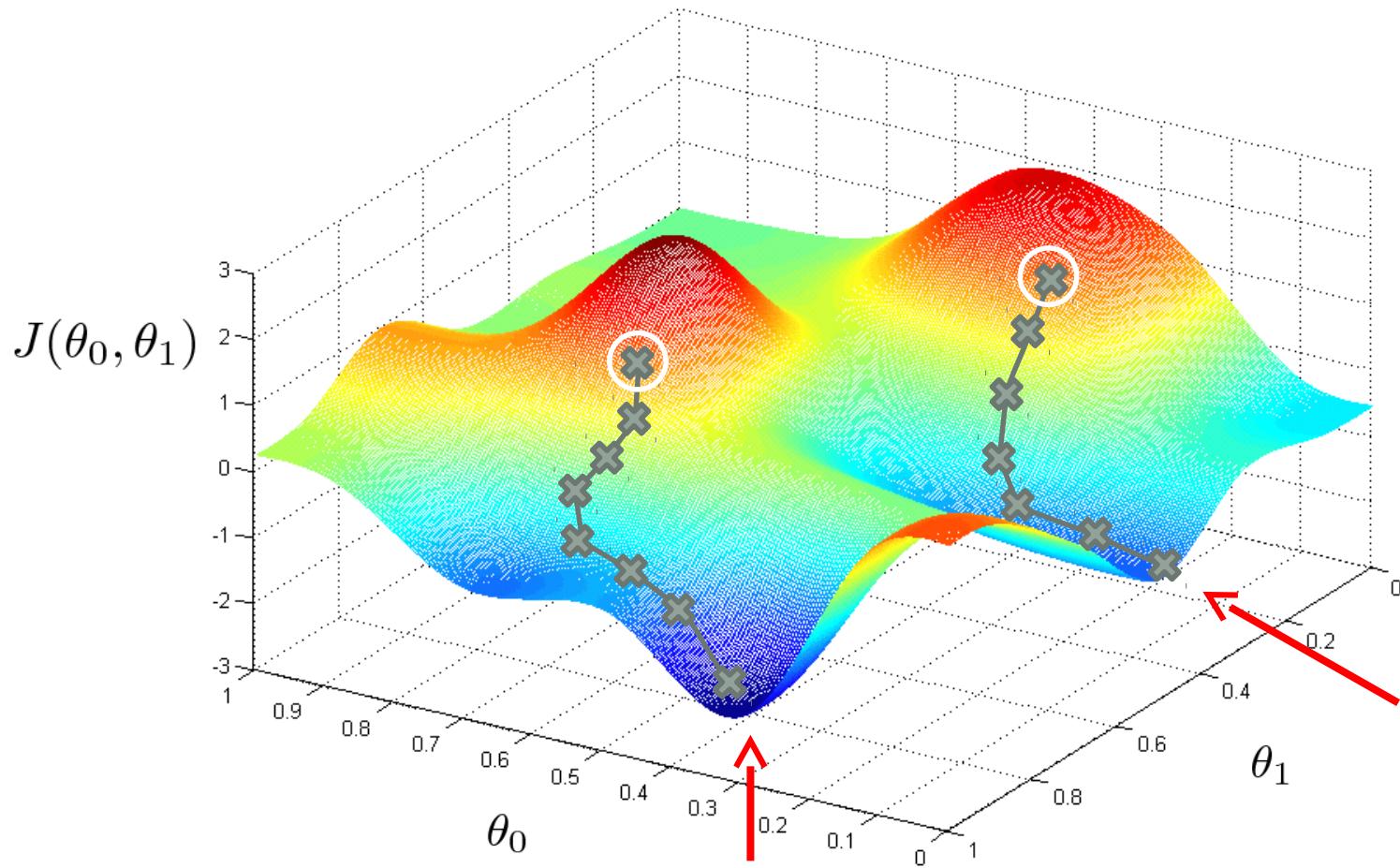
$$J(\theta_0, \theta_1)$$



EXCURSUS: GRADIENT DESCENT

Descending in the steepest direction

Gradient descent on some arbitrary cost function $J(\theta_0, \theta_1) \dots$



Gradient descent algorithm

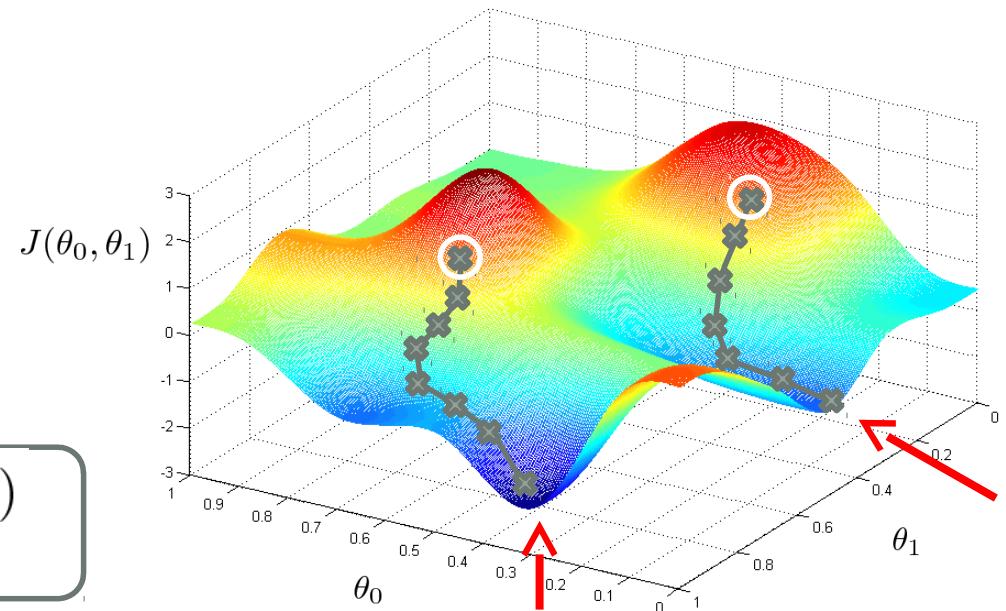
- Repeat until convergence

$$\theta_j := \theta_j - \eta \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{simultaneously updating } \theta_0 \text{ and } \theta_1)$$

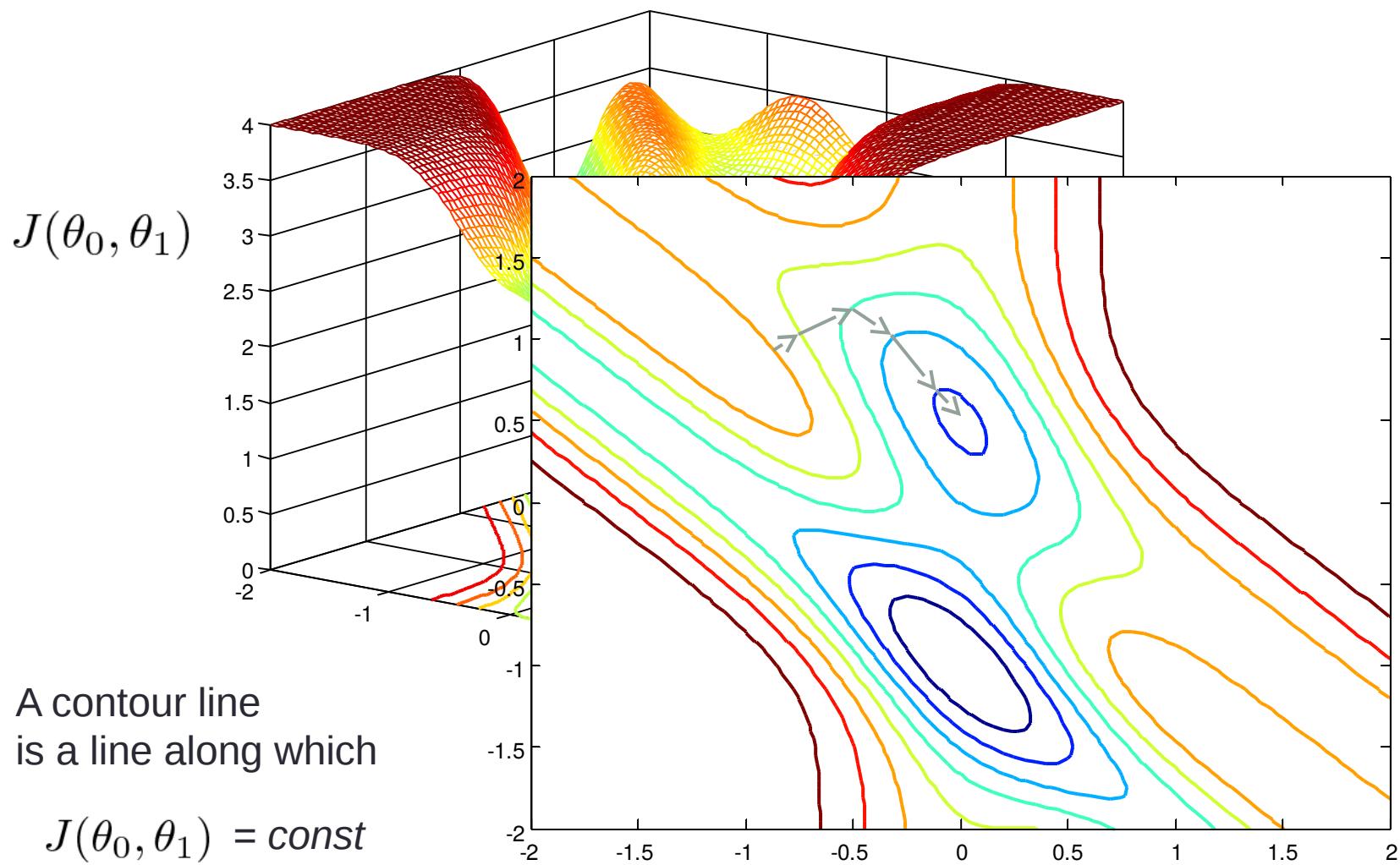
negative gradient =
descent

learning rate („eta“)

partial derivative of
 $J(\theta_0, \theta_1)$
with respect to θ_j

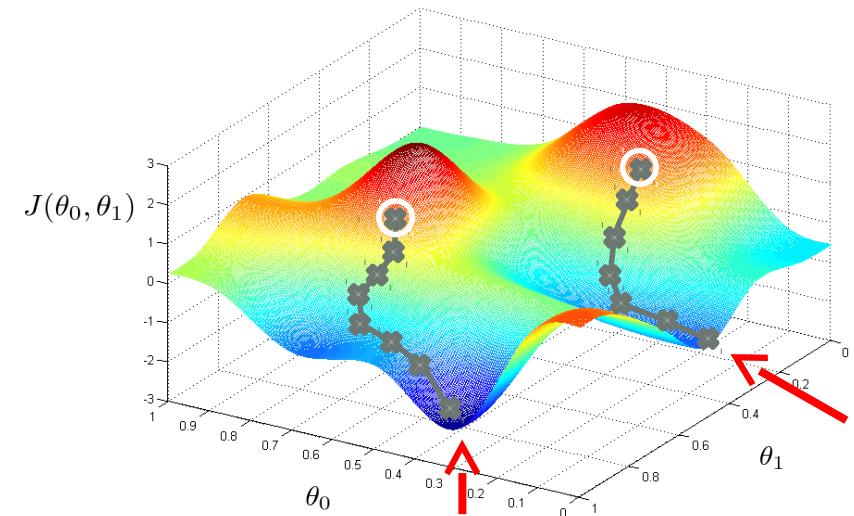


Gradient is orthogonal to contour lines

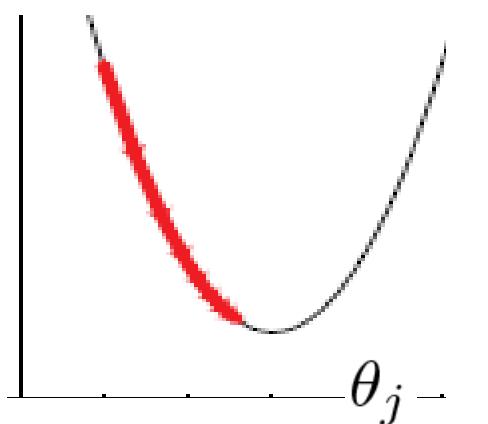


Potential issues with gradient descent

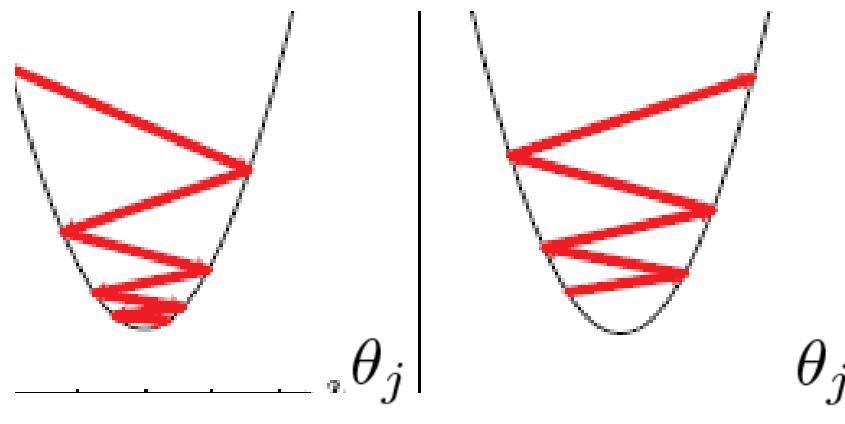
- May get stuck in local minima
- Learning rate too small: slow convergence
- Learning rate too large: oscillations, divergence



η too small



η too large



LINEAR REGRESSION WITH GRADIENT DESCENT

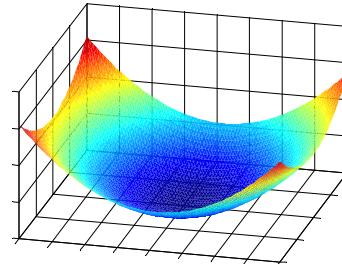
(ONE INPUT)

Application of gradient descent

- Linear regression cost

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$



- Gradient descent

$$\theta_j := \theta_j - \eta \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(simultaneous update)

$$\theta_0 := \theta_0 - 2\eta \cdot \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)$$

$$\theta_1 := \theta_1 - 2\eta \cdot \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right) \cdot x^{(i)}$$

(simultaneous update)

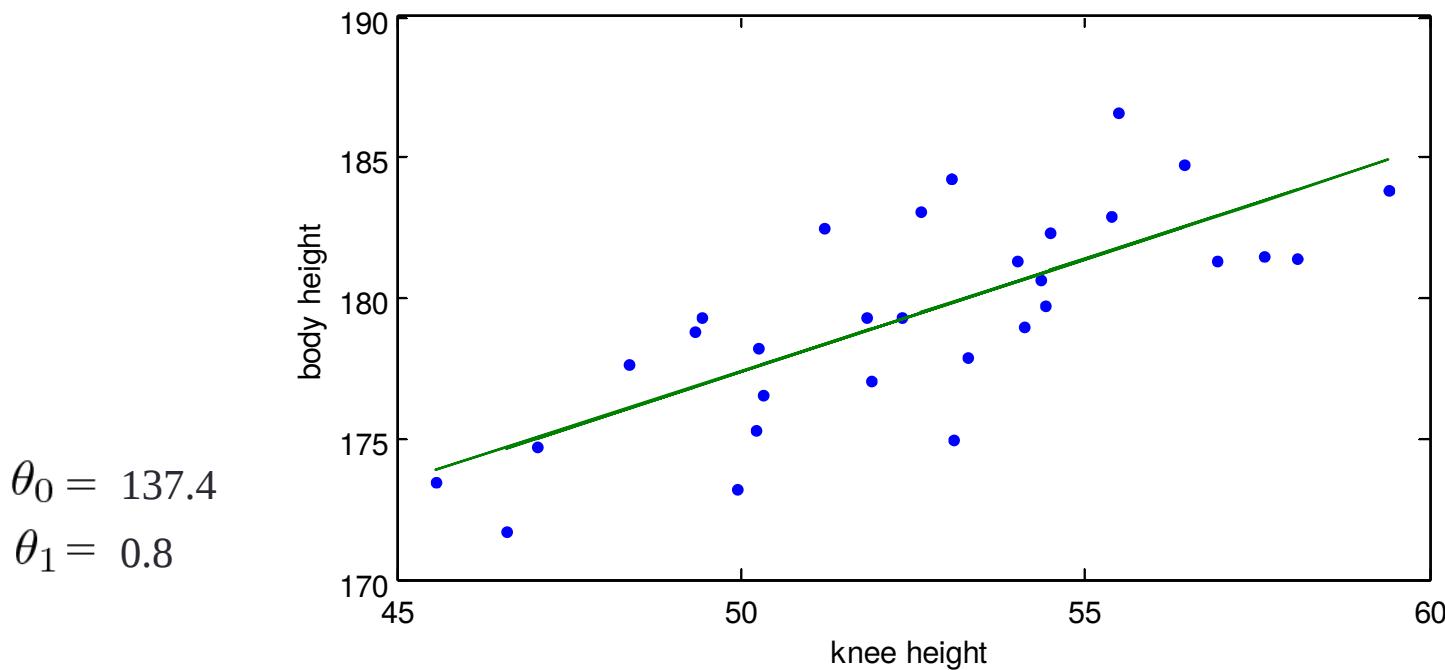
learning rate

„error“

„input“

Predicting height from knee height

- Optimal fit to training data



LINEAR REGRESSION

MORE GENERAL FORMULATION: MULTIPLE FEATURES

Multiple inputs (features)

	Knee Height x_1	Arm span x_2	Age x_3	Height y
	50	166	32	171
m	56	172	17	175
	52	174	62	168

$n = 3$

- Notation:

m ... number of training examples

n ... number of features

$\boldsymbol{x}^{(i)}$... input features of i 'th training example (vector-valued)
 $x_j^{(i)}$... value of feature j in i 'th training example

$$\boldsymbol{x}^{(2)} = \begin{pmatrix} 56 \\ 172 \\ 17 \end{pmatrix}$$

$$x_3^{(2)} = 17$$

Linear hypothesis

- Hypothesis (one input):

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

- Hypothesis (multiple input features):

$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 \cdot x_1 + \cdots + \theta_n \cdot x_n$$

Example: $h(x) = 50 + 0.5 \cdot \text{kneeheight} + 0.3 \cdot \text{armspan} + 0.1 \cdot \text{age}$

- More compact notation:

$$h_{\theta}(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta}$$

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \boldsymbol{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix}$$

Introduce $x_0 = 1$
Why? Notation convenience!

Multiple inputs (features) revisited

	x_0	Knee Height x_1	Arm span x_2	Age x_3	Height y
	1	50	166	32	171
m	1	56	172	17	175
	1	52	174	62	168
	1

$n = 3$

- Notation:

m ... number of training examples

n ... number of features

$\boldsymbol{x}^{(i)}$... input features of i 'th training example (vector-valued)

$x_j^{(i)}$... value of feature j in i 'th training example

$$\boldsymbol{x}^{(2)} = \begin{pmatrix} 1 \\ 56 \\ 172 \\ 17 \end{pmatrix}$$

$$x_0^{(2)} = 1$$

$$x_3^{(2)} = 17$$

Matrix and vector notation

x_0	Knee Height x_1	Arm span x_2	Age x_3	Height y
1	50	166	32	171
1	56	172	17	175
1	52	174	62	168

$$\mathbf{X} = \begin{pmatrix} 1 & 50 & 166 & 32 \\ 1 & 56 & 172 & 17 \\ 1 & 52 & 174 & 62 \end{pmatrix}$$

$$\mathbf{y} = \begin{pmatrix} 171 \\ 175 \\ 168 \end{pmatrix}$$

$$\mathbf{x}^{(i)} = \begin{pmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} \quad (\mathbf{x}^{(1)})^T \quad \\ \quad (\mathbf{x}^{(2)})^T \quad \\ \vdots \\ \quad (\mathbf{x}^{(m)})^T \quad \end{pmatrix}$$

features of i'th training example
 $(n+1) \times 1$

design matrix
 $m \times (n+1)$

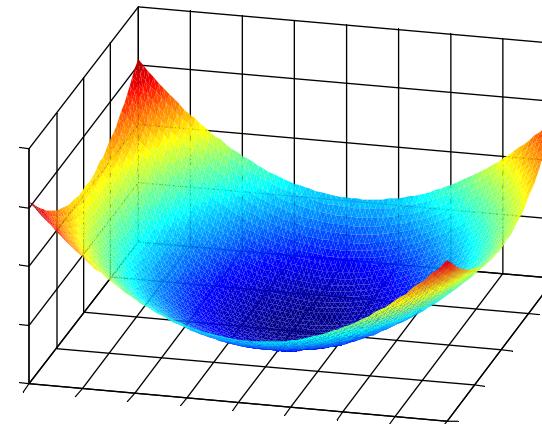
output/target vector
 $m \times 1$

LINEAR REGRESSION WITH GRADIENT DESCENT (GENERAL FORMULATION)

Linear regression problem statement

- Hypothesis: $h_{\theta}(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta}$
- Cost function: $J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$

high-dimensional quadratic
("bowl"-shaped) function



Goal is to find parameters which minimize the cost

Gradient descent (multiple features)

with **one** input feature:

$$\theta_0 := \theta_0 - 2\eta \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

learning rate

$$\theta_1 := \theta_1 - 2\eta \cdot \frac{1}{m} \sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})}_{\text{„error“}} \cdot \underbrace{x^{(i)}}_{\text{„input“}}$$

(simultaneous update)

with **n** input features:

$$\theta_j := \theta_j - 2\eta \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

learning rate

(simultaneous update for $j=0 \dots n$)

For $j=0$: define for convenience $x_0^{(i)} = 1$

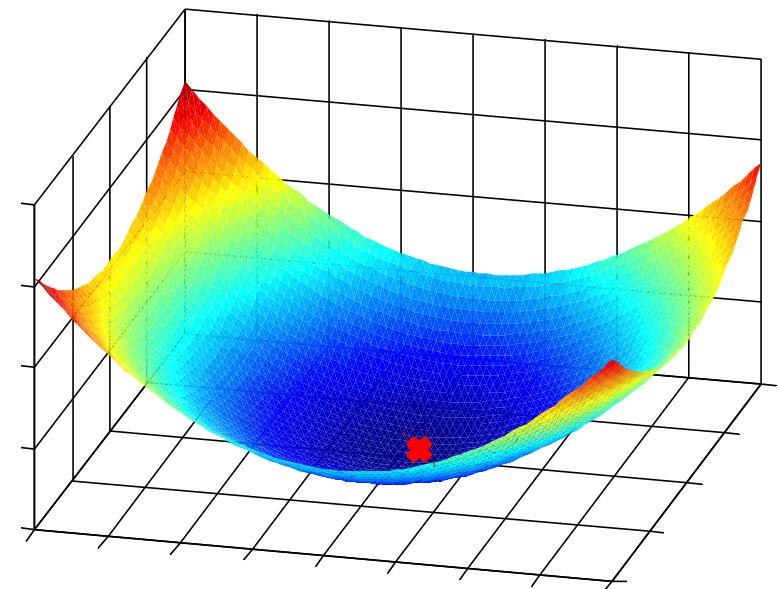
LINEAR REGRESSION ANALYTICAL SOLUTION

Analytical solution

- Set all partial derivatives of cost function $J(\theta) = 0$
- Solving system of linear equations yields:

$$\boxed{\theta^* = \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}}$$

Moore-Penrose Pseudoinverse of \mathbf{X}



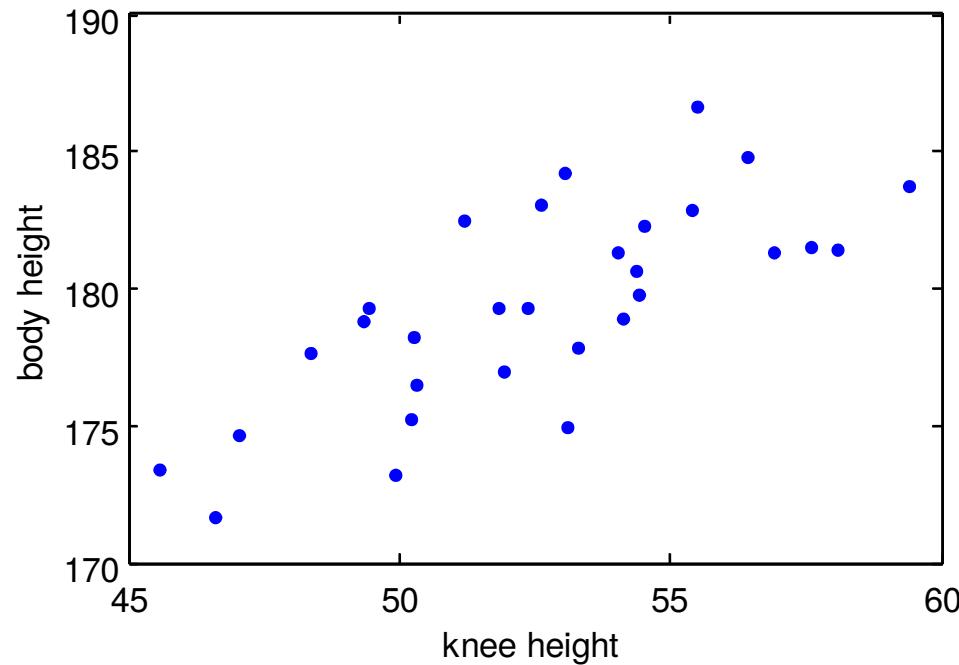
\mathbf{X} ... design matrix

\mathbf{y} ... output/target vector

- Note: This analytical solution requires that columns of \mathbf{X} are linearly independent („regular“ conditions)

Example: analytical solution applied to problem with one input

Knee Height [cm]	Height [cm]
50	171
56	175
52	168
...	...



Example: analytical solution applied to problem with one input

Knee Height [cm]	Height [cm]
50	171
56	175
52	168
...	...

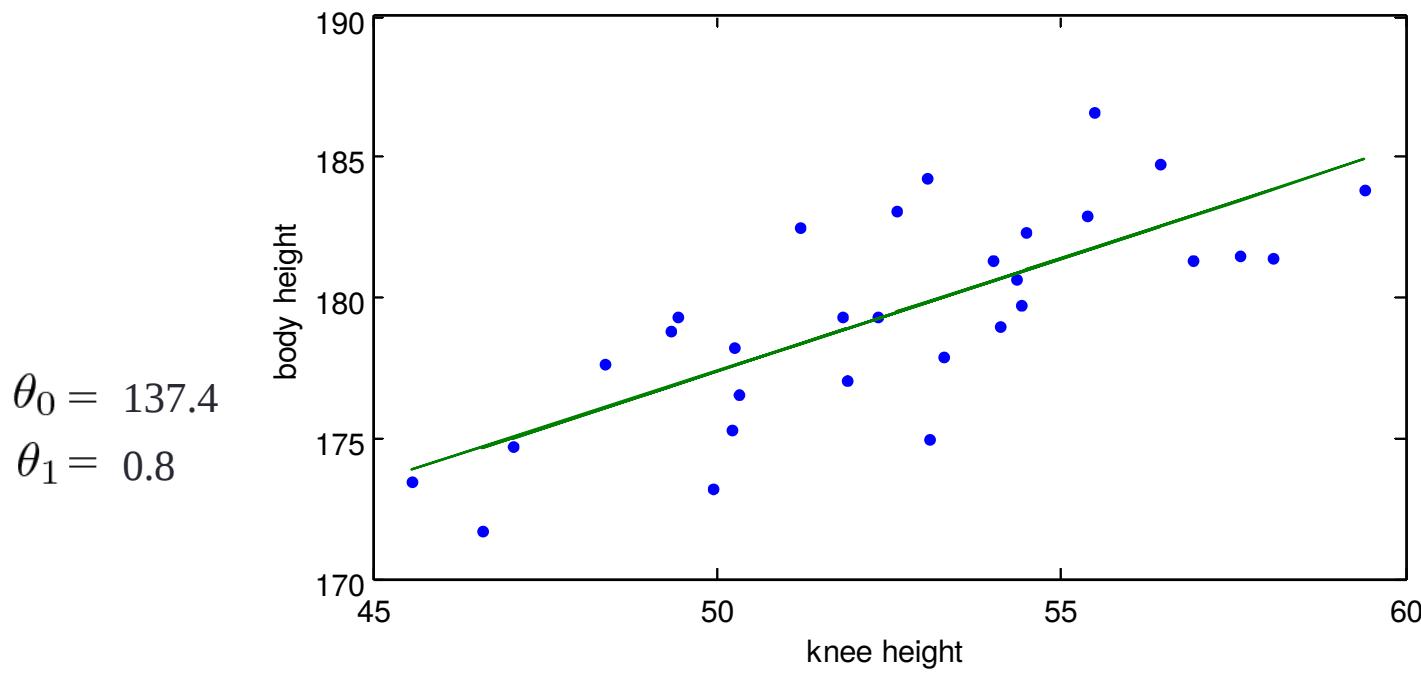
$$\mathbf{X} = \begin{pmatrix} 1 & 50 \\ 1 & 56 \\ 1 & 52 \\ \vdots & \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 171 \\ 175 \\ 168 \\ \vdots \end{pmatrix}$$

30×2 30×1

$$\begin{aligned}\boldsymbol{\theta}^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \begin{pmatrix} 137.4 \\ 0.8 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{X}^T \mathbf{X} &= \begin{pmatrix} 30 & 1577 \\ 1577 & 83222 \end{pmatrix} & 2 \times 2 \\ (\mathbf{X}^T \mathbf{X})^{-1} &= \begin{pmatrix} 7.994 & -0.152 \\ -0.152 & 0.003 \end{pmatrix} & 2 \times 2 \\ \mathbf{X}^T \mathbf{y} &= \begin{pmatrix} 5383 \\ 283210 \end{pmatrix} & 2 \times 1\end{aligned}$$

Predicting height from knee height



$$\begin{aligned}\boldsymbol{\theta}^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \begin{pmatrix} 137.4 \\ 0.8 \end{pmatrix}\end{aligned}$$

Gradient descent

- Need to choose learning rate η
- Iterative algorithm (needs many iterations to converge)
- Works well even when number of input features n is large

Analytical solution

- No need to choose η
- Direct solution (no iteration)
- Slow if n is too large (inverting $n \times n$ matrix)

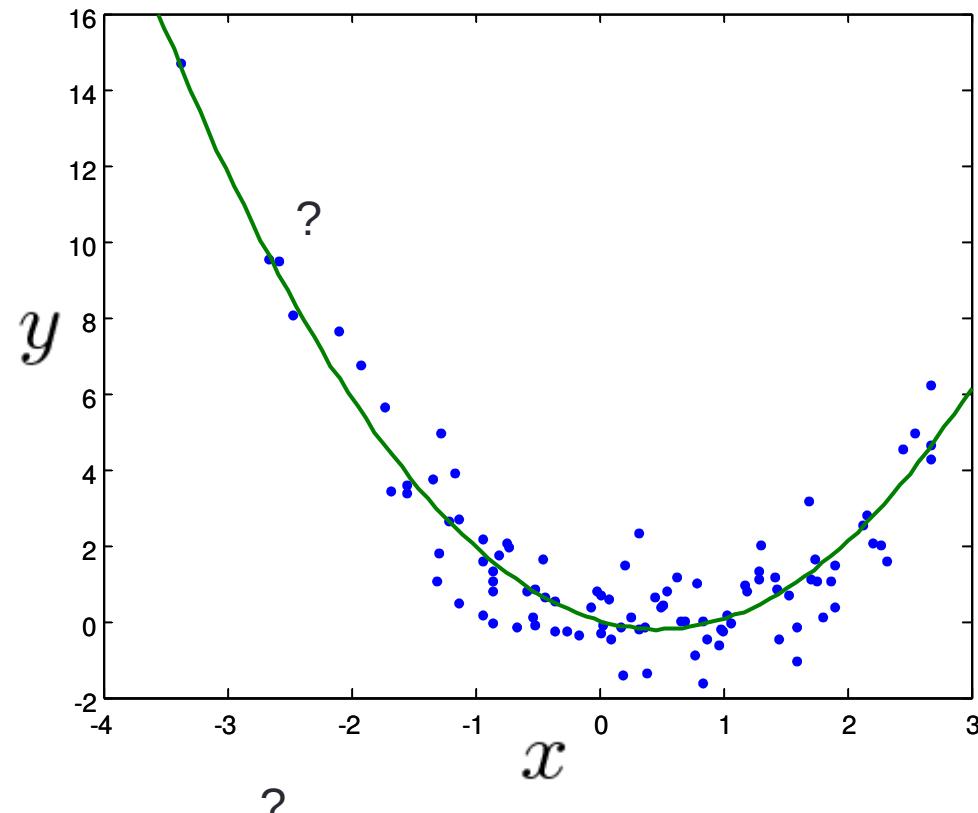
NON-LINEAR FEATURES

(NON-LINEAR BASIS FUNCTIONS)

Non-linear trends in data

- How can we learn non-linear hypotheses?

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...



$$h_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 \cdot x + \theta_2 \cdot x^2$$

Linear fit to this „non-linear“ data

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...

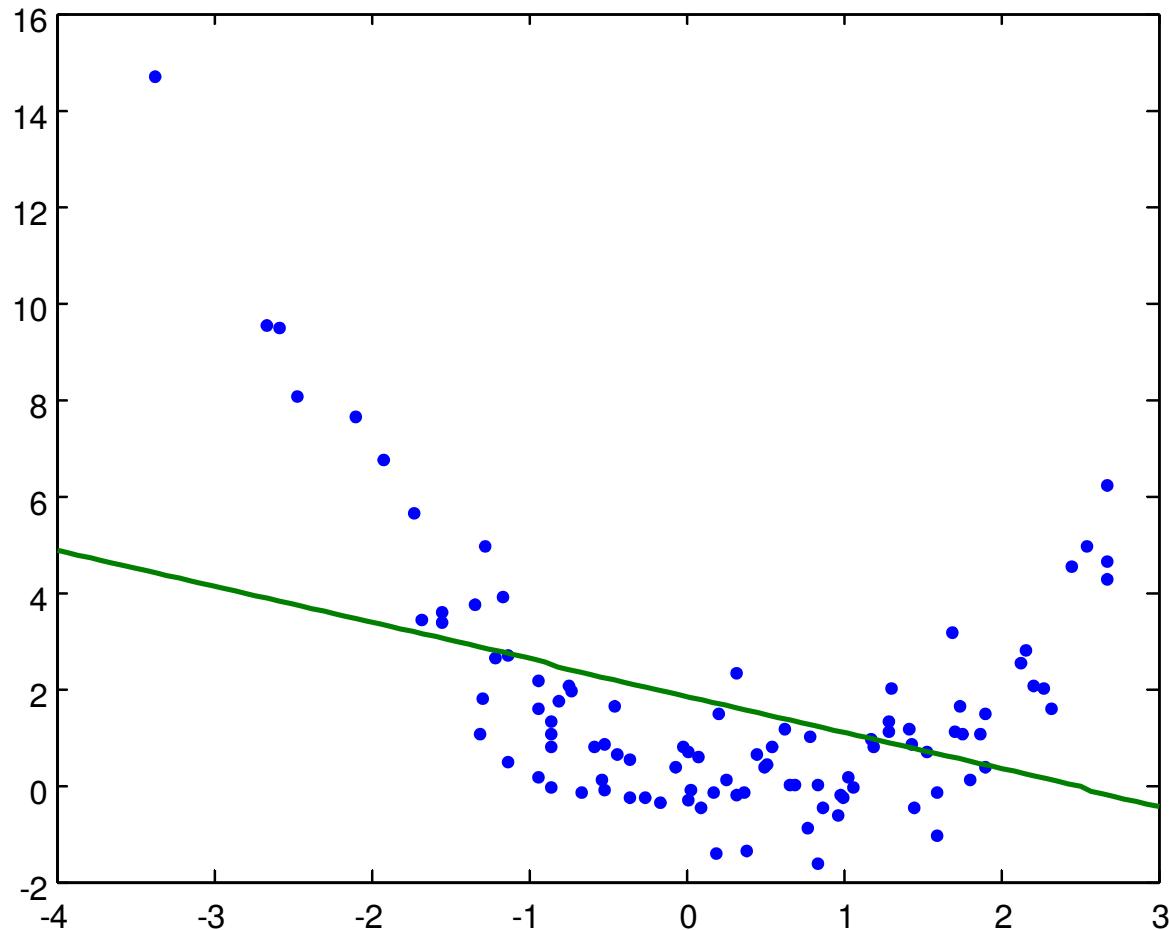
$$\mathbf{X} = \begin{pmatrix} 1 & 0.01 \\ 1 & -1.22 \\ 1 & 0.17 \\ \vdots \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} -0.27 \\ 2.63 \\ -0.13 \\ \vdots \end{pmatrix}$$

standard design matrix

Hypothesis: $h_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 \cdot x$

Optimal parameters: $\boldsymbol{\theta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Linear fit to this „non-linear“ data



$$h_{\theta}(x) = 1.85 - 0.76 \cdot x$$

Non-linear (quadratic) fit

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...

$$\phi_0 = 1 \quad \phi_1 = x \quad \phi_2 = x^2$$

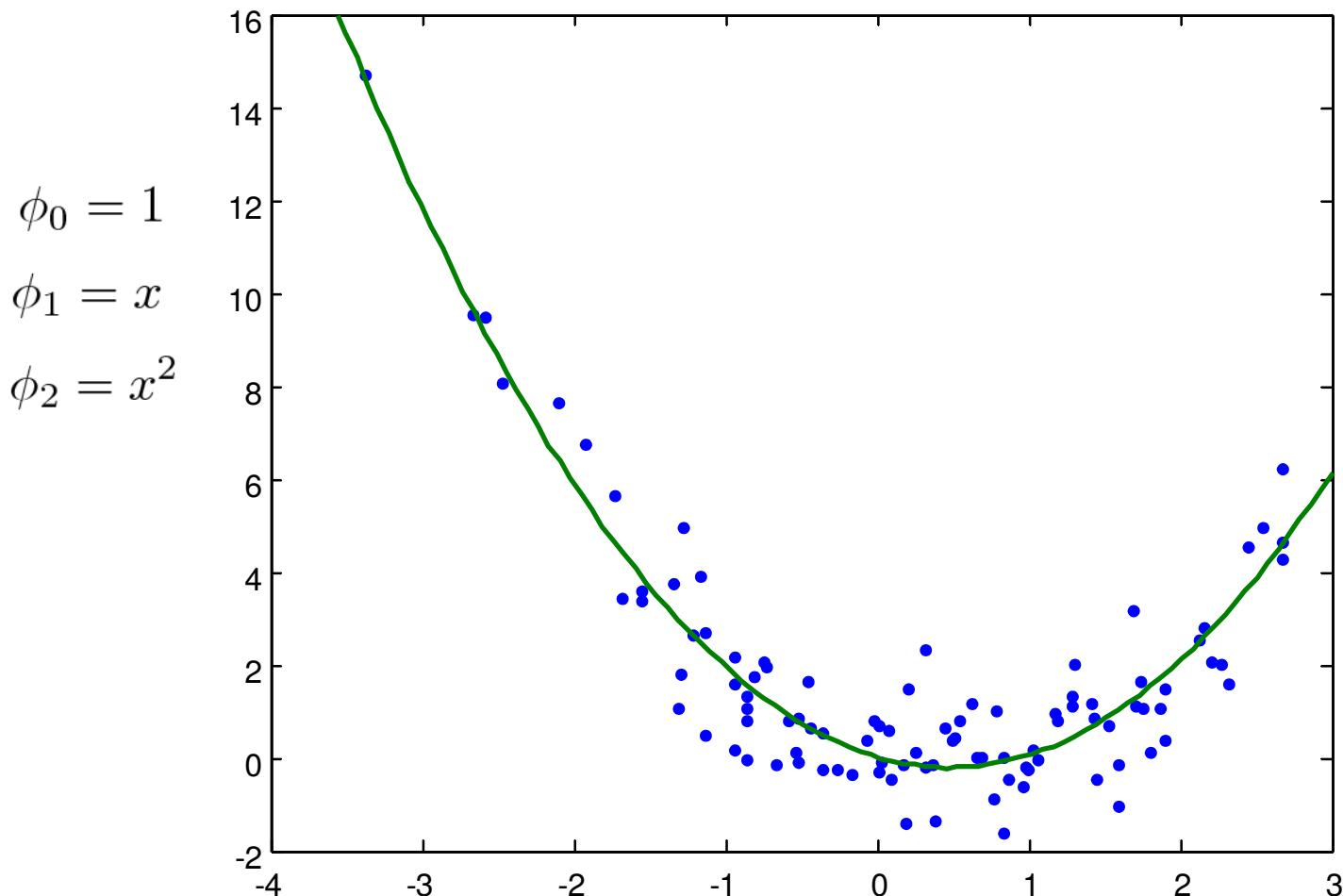
$$\Phi = \begin{pmatrix} 1 & 0.01 & 0.01^2 \\ 1 & -1.22 & (-1.22)^2 \\ 1 & 0.17 & (0.17)^2 \\ \vdots & & \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} -0.27 \\ 2.63 \\ -0.13 \\ \vdots \end{pmatrix}$$

*design matrix with
non-linear features*

Hypothesis: $h_{\boldsymbol{\theta}}(\boldsymbol{\phi}) = \theta_0 + \theta_1 \cdot \phi_1 + \theta_2 \cdot \phi_2$

Optimal parameters: $\boldsymbol{\theta}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$

Non-linear (quadratic) fit



$$h_{\theta}(x) = 0.02 \cdot 1 - 0.95 \cdot x + 0.99 \cdot x^2$$

Non-linear (sinusoid) fit

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...

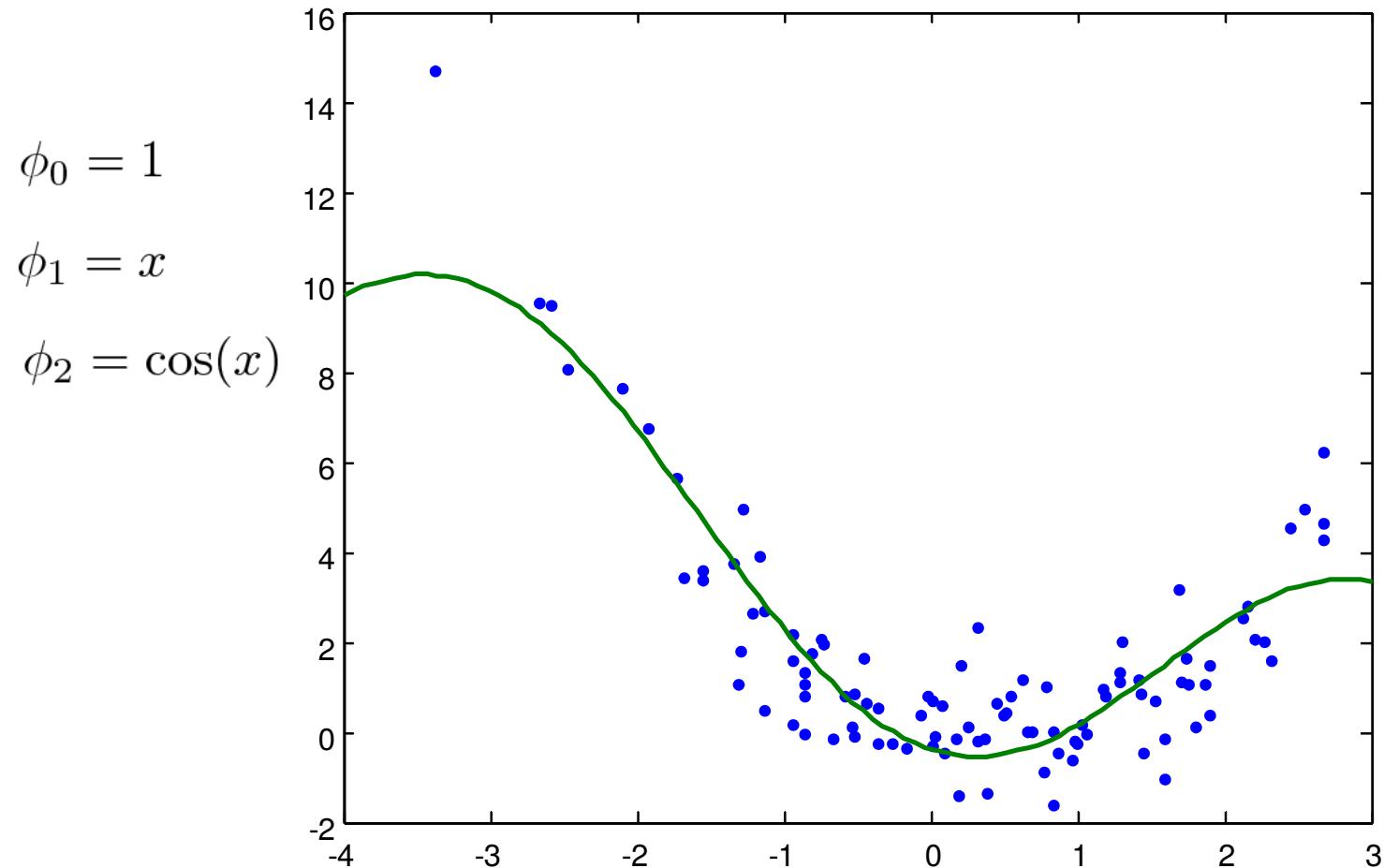
$$\Phi = \begin{pmatrix} 1 & 0.01 & \cos(0.01) \\ 1 & -1.22 & \cos(-1.22) \\ 1 & 0.17 & \cos(0.17) \\ \vdots & & \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} -0.27 \\ 2.63 \\ -0.13 \\ \vdots \end{pmatrix}$$

*design matrix with
non-linear features*

Hypothesis: $h_{\boldsymbol{\theta}}(\boldsymbol{\phi}) = \theta_0 + \theta_1 \cdot \phi_1 + \theta_2 \cdot \phi_2$

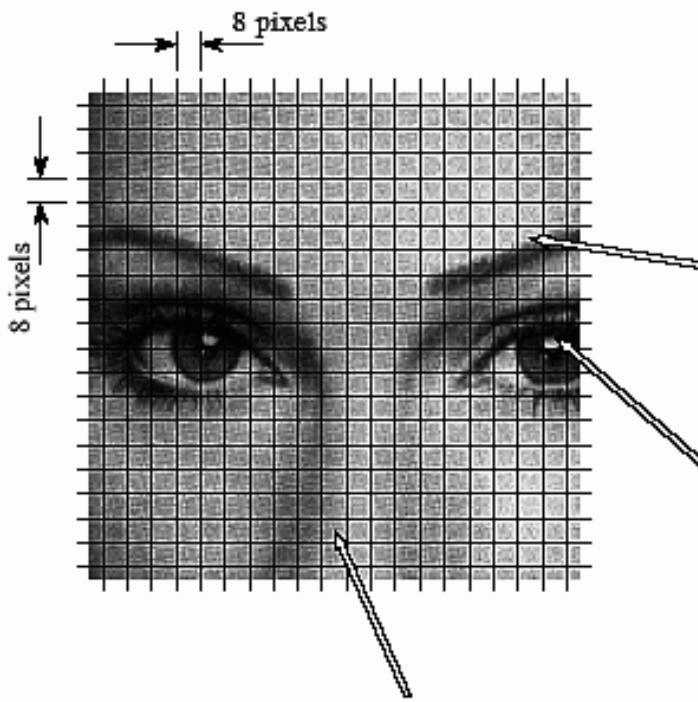
Optimal parameters: $\boldsymbol{\theta}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$

Non-linear (sinusoid) fit

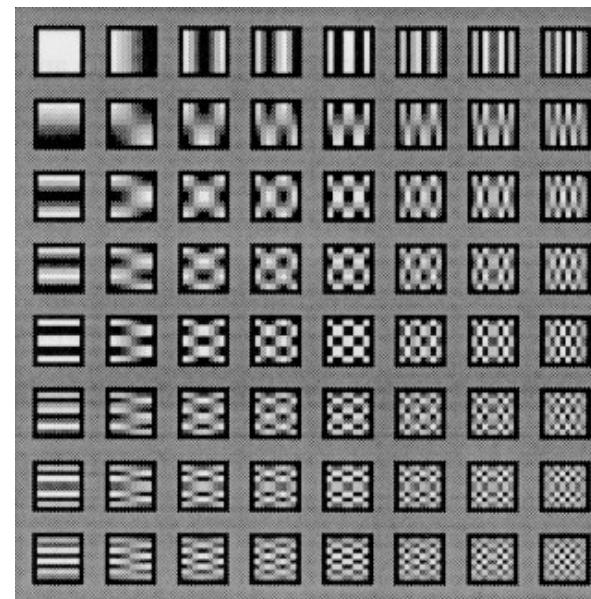


$$h_{\theta}(x) = 3.12 \cdot 1 - 1.07 \cdot x - 3.5 \cdot \cos(x)$$

Image: JPEG = cosin-basis



Each block of 8x8 pixels is represented in a Fourier basis of **cosin filters**

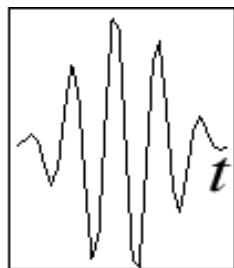


Better representation of **edges and corners**
Allows for compression

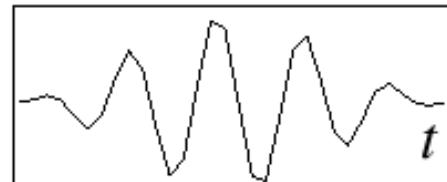
Audio: cosin or wavelet basis

Good signal representation make a compromise between time and frequency

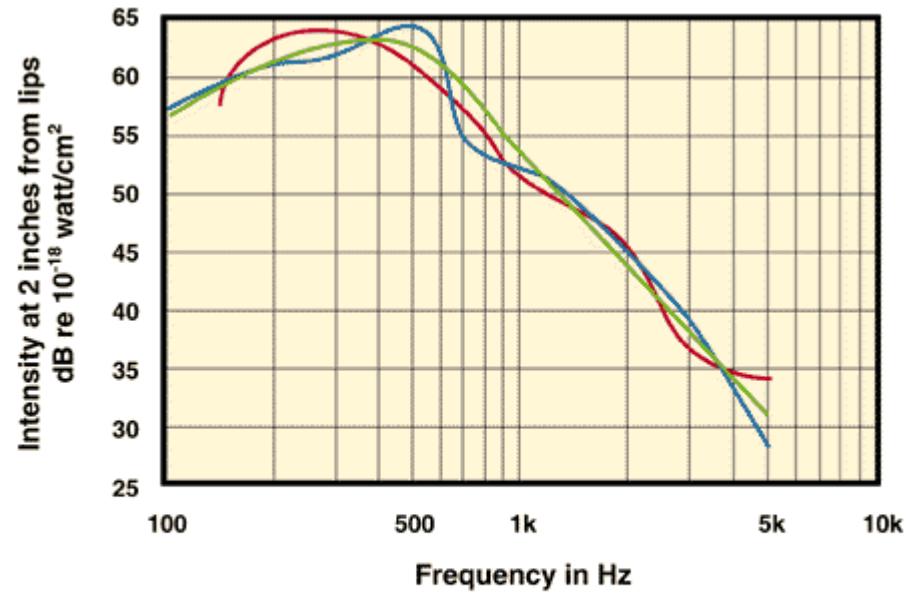
$f(t)$
(Mother
Wavelet)



$\frac{1}{\sqrt{2}}f\left(\frac{t}{2}\right)$



$\frac{1}{2}f\left(\frac{t}{4}\right)$



- Avg. of five women
- Avg. of five men
- Idealized Vocal Spectrum

Non-linear input features (in general)

feature 2 of all training examples

$$\Phi = \begin{pmatrix} 1 & \phi_1^{(1)} & \boxed{\phi_2^{(1)}} & \dots & \phi_n^{(1)} \\ 1 & \phi_1^{(2)} & \phi_2^{(2)} & \dots & \phi_n^{(2)} \\ \vdots & & & & \\ 1 & \phi_1^{(m)} & \phi_2^{(m)} & \dots & \phi_n^{(m)} \end{pmatrix}$$

*all features of
1st training example*

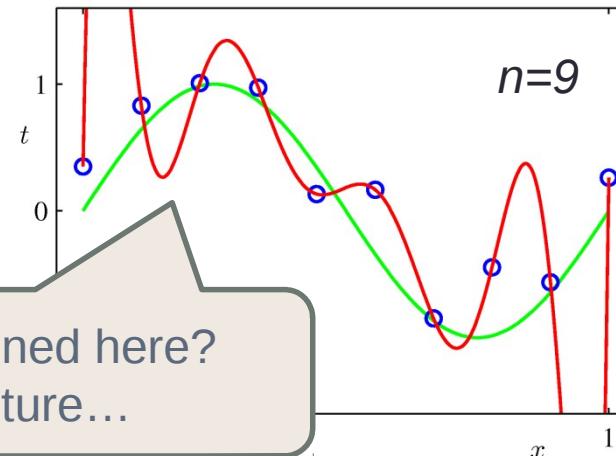
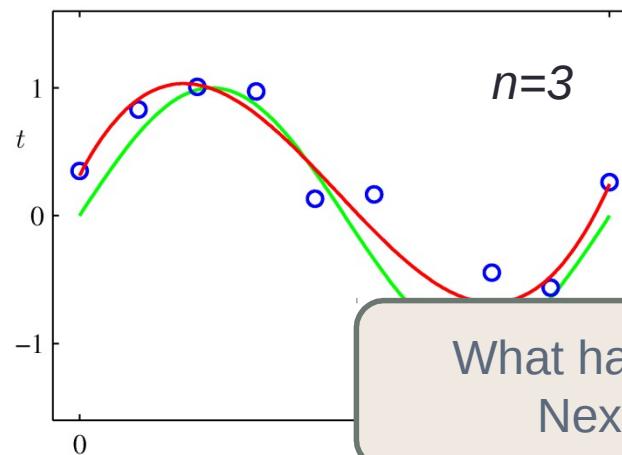
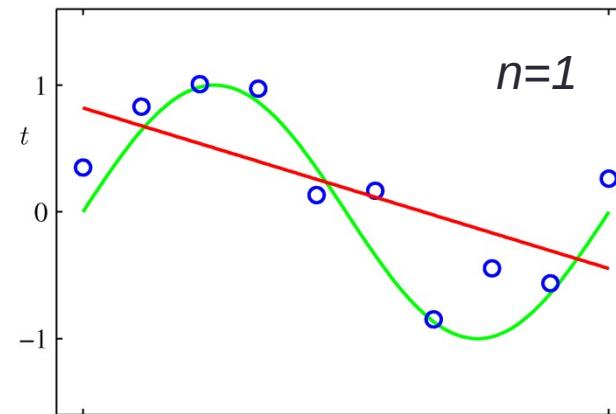
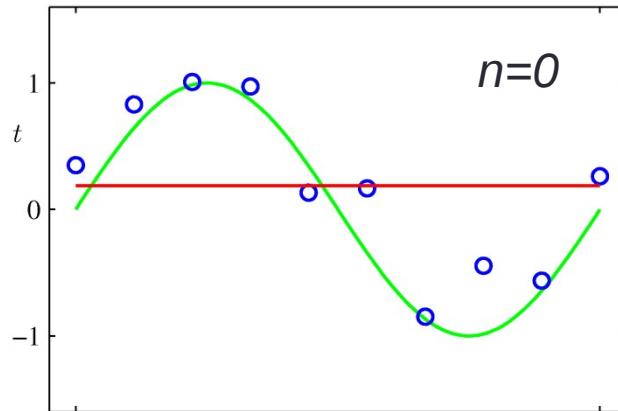
- Feature 2 for each training example i is computed by applying a **non-linear basis function**: $\phi_2^{(i)} = \phi_2(\mathbf{x}^{(i)})$
- Allows to learn a variety of non-linear functions with the same technique(s):
 - Gradient descent or $\theta^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$

Polynomial regression

- Features are powers of x

$$\phi_0 = x^0, \phi_1 = x^1, \phi_2 = x^2, \dots, \phi_n = x^n$$

$n =$ degree of polynome
to be learned

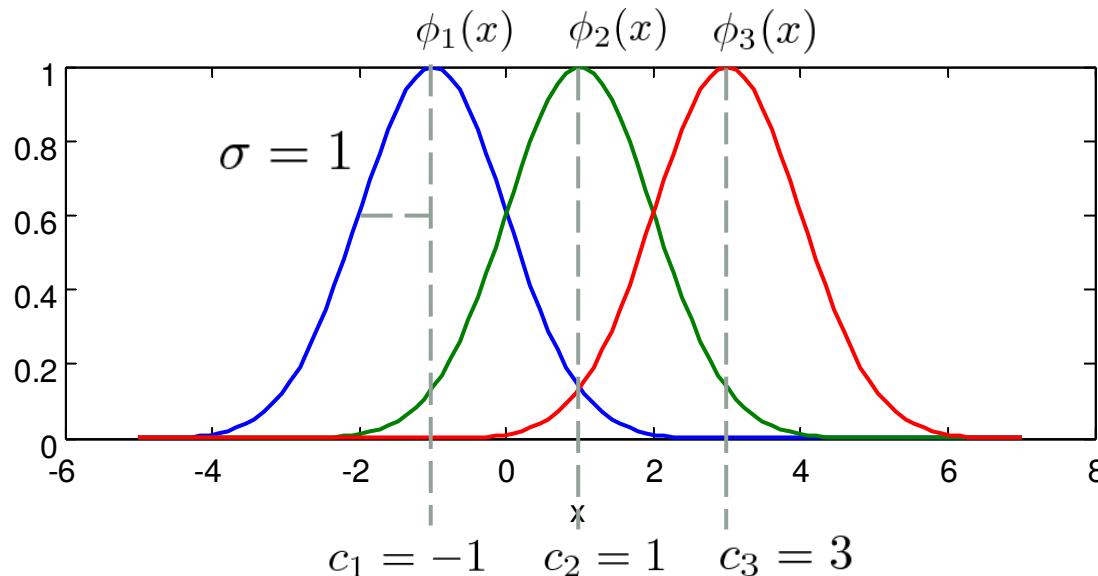


What happened here?
Next lecture...

Radial basis functions

- „Gaussian“-shaped RBFs:
 - Each basis function j has a **center** \mathbf{c}_j in the input space
 - The **width** of the basis functions is determined by σ

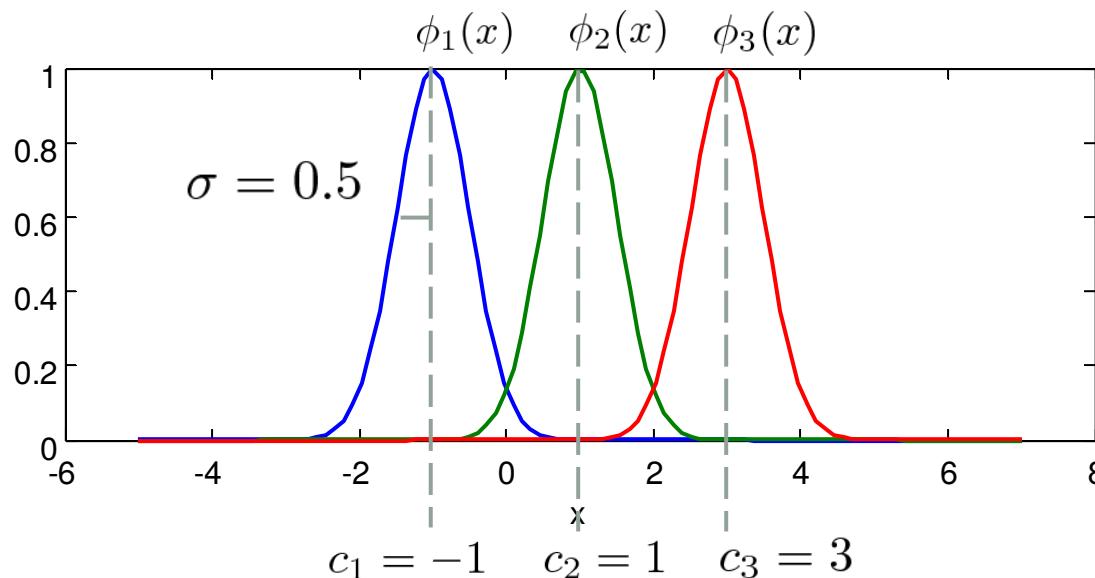
$$\phi_j(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \cdot \|\mathbf{x} - \mathbf{c}_j\|^2\right)$$



Radial basis functions

- „Gaussian“-shaped RBFs:
 - Each basis function j has a **center** \mathbf{c}_j in the input space
 - The **width** of the basis functions is determined by σ

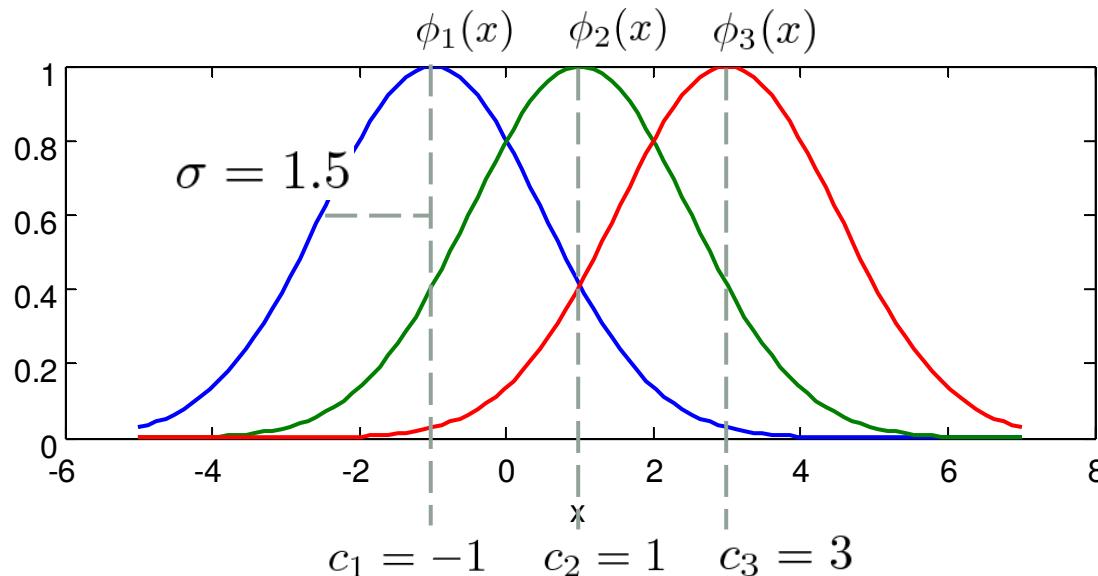
$$\phi_j(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \cdot \|\mathbf{x} - \mathbf{c}_j\|^2\right)$$



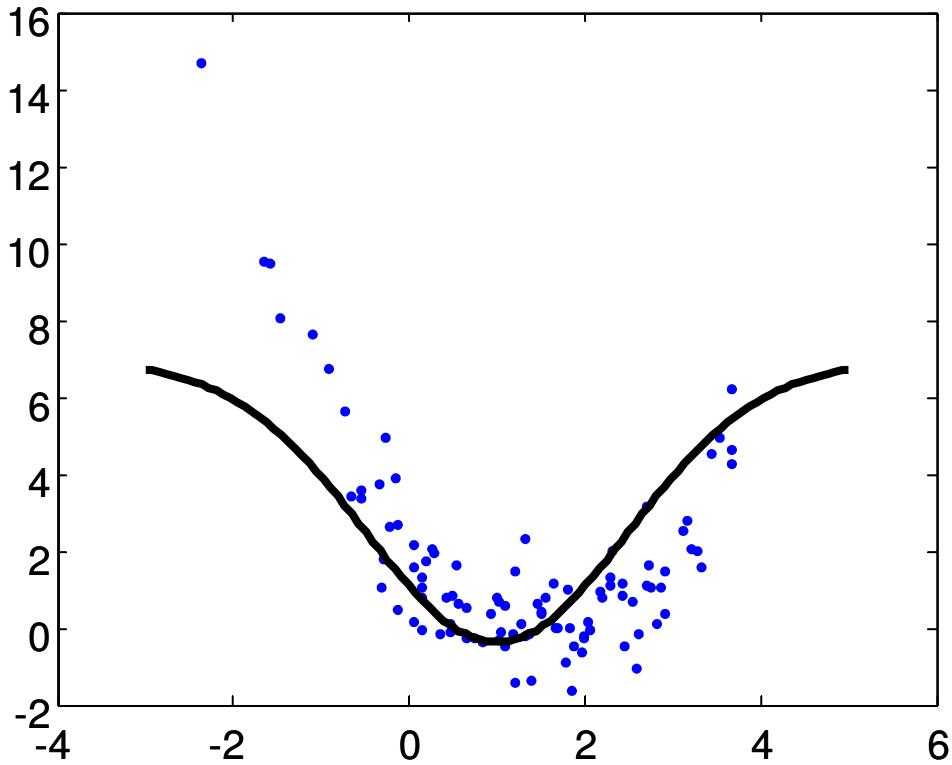
Radial basis functions

- „Gaussian“-shaped RBFs:
 - Each basis function j has a **center** \mathbf{c}_j in the input space
 - The **width** of the basis functions is determined by σ

$$\phi_j(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \cdot \|\mathbf{x} - \mathbf{c}_j\|^2\right)$$



Fitting a single RBF to data

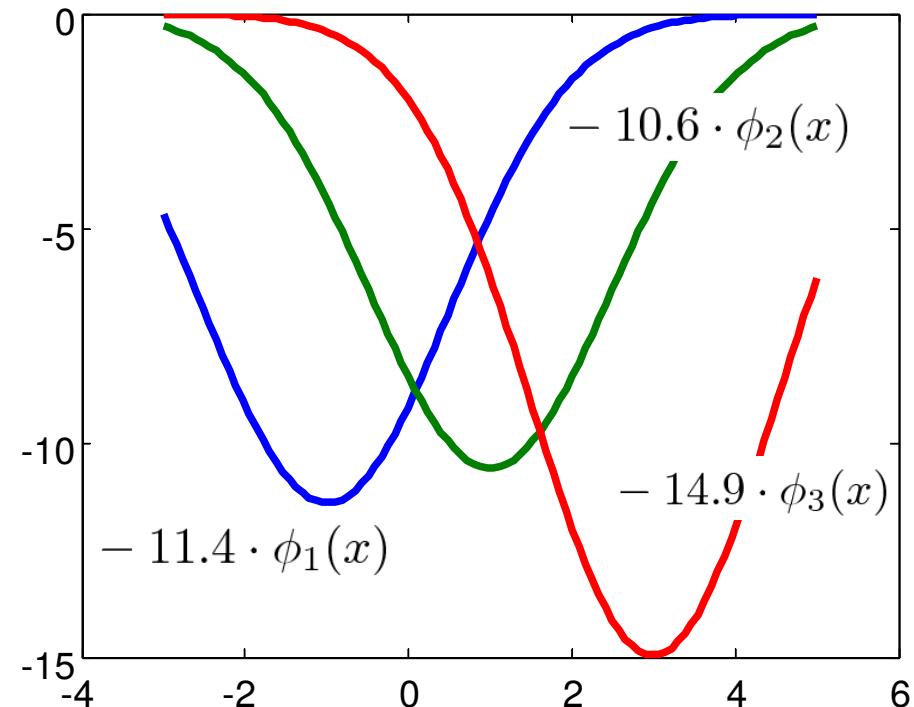
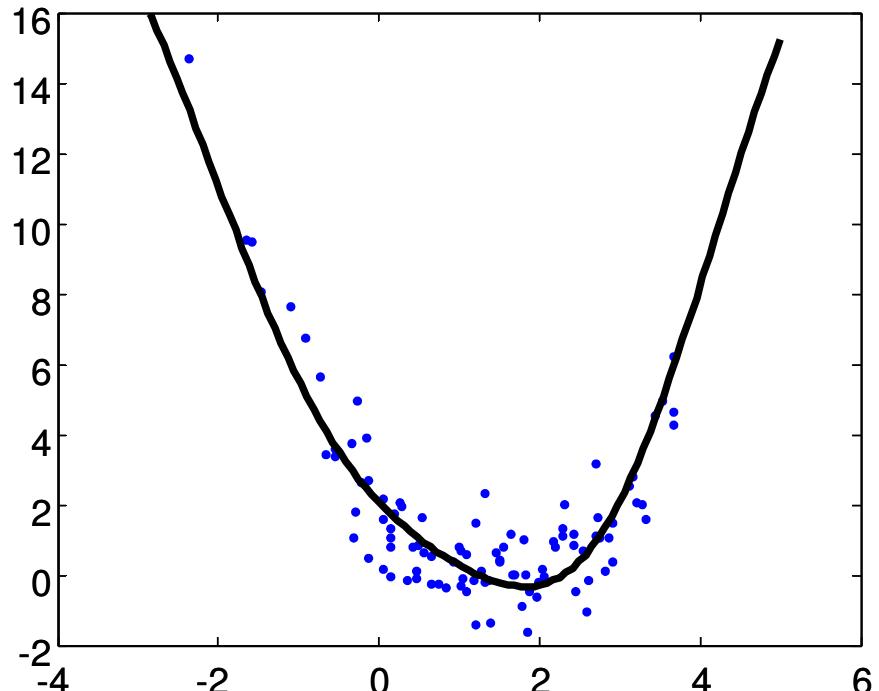


$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \phi_1(x)$$

$$h_{\theta}(x) = 6.9 - 7.3 \cdot \phi_1(x)$$

RBF with $c_1 = 1$
 $\sigma = 1.5$

Fitting RBFs to data



RBFs with $\sigma = 1.5$

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \phi_1(x) + \theta_2 \cdot \phi_2(x) + \theta_3 \cdot \phi_3(x)$$

$$h_{\theta}(x) = 21.7 - 11.4 \cdot \phi_1(x) - 10.6 \cdot \phi_2(x) - 14.9 \cdot \phi_3(x)$$

SUMMARY (QUESTIONS)

Some questions...

- Hypothesis for linear regression = ?
- Cost function for linear regression = ?
- How many local minima may the cost function for lin. reg. have (under regular conditions)?
- Name two ways to minimize the cost function?
- General gradient descent formula?
- Linear regression with gradient descent formula?
- What issues can arise during gradient descent?
- What is the design matrix? What are its dimensions?
- Analytical solution for linear regression = ?
 - What are the components of the solution?
- Pros and Cons of gradient descent vs. analytical solution?
- How can one learn non-linear hypotheses with linear regression?
- What is polynomial regression?
- What are radial basis functions?

What is next?

- Classification with Logistic Regression
- Gradient descent tricks & more advanced optimization techniques
- Underfitting & Overfitting
- Model selection (Training- & Validation- & Testset)